(12) **United States Patent**
Megiddo et al.

(10) **Patent No.:** **US 6,996,676 B2**
(45) **Date of Patent:** **Feb. 7, 2006**

(54) **SYSTEM AND METHOD FOR IMPLEMENTING AN ADAPTIVE REPLACEMENT CACHE POLICY**

(75) Inventors: **Nimrod Megiddo**, Palo Alto, CA (US); **Dharmendra Shantilal Modha**, San Jose, CA (US)

(73) Assignee: **International Business Machines Corporation**, Armonk, NY (US)

( * ) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 465 days.

(21) Appl. No.: **10/295,507**

(22) Filed: **Nov. 14, 2002**

(65) **Prior Publication Data**

US 2004/0098541 A1 May 20, 2004

(51) **Int. Cl.**
*G06F 12/00* (2006.01)

(52) **U.S. Cl.** ........................ **711/129**; 711/133; 711/170
(58) **Field of Classification Search** ................. 711/129, 711/133, 170
See application file for complete search history.

(56) **References Cited**

U.S. PATENT DOCUMENTS

| | | | | | |
|---|---|---|---|---|---|
| 4,463,424 | A | * | 7/1984 | Mattson et al. | .............. 711/136 |
| 4,464,712 | A | | 8/1984 | Fletcher | ....................... 364/200 |
| 4,503,501 | A | * | 3/1985 | Coulson et al. | ............. 711/129 |
| 4,780,815 | A | * | 10/1988 | Shiota | ......................... 711/171 |
| 5,481,691 | A | | 1/1996 | Day, III et al. | ............. 395/425 |
| 5,752,255 | A | * | 5/1998 | Jarvis | ............................ 711/3 |
| 6,041,390 | A | | 3/2000 | Liu et al. | ..................... 711/110 |
| 6,154,813 | A | | 11/2000 | Martin et al. | ................ 711/133 |

| | | | | |
|---|---|---|---|---|
| 6,209,062 | B1 | 3/2001 | Boland et al. | .............. 711/134 |
| 6,327,643 | B1 | 12/2001 | Egan | ........................... 711/134 |
| 6,408,368 | B1 | 6/2002 | Parady | ....................... 711/159 |

OTHER PUBLICATIONS

"Least–Recently–Used–Page–Replacement Algorithm For Cache Memories," IBM Technical Disclosure Bulletin, vol. 25 No. 3A, Aug. 1982.
S. Kim et al., "Area Efficient Architectures for Information Integrity in Cache Memories," IEEE–CS\TCCA:TC on Computer Architecture, 1999.
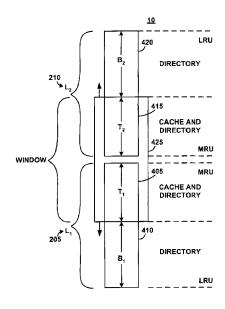
(Continued)

*Primary Examiner*—Jack Lane
(74) *Attorney, Agent, or Firm*—Samuel A. Kassatly

(57) **ABSTRACT**

An adaptive replacement cache policy dynamically maintains two lists of pages, a recency list and a frequency list, in addition to a cache directory. The policy keeps these two lists to roughly the same size, the cache size c. Together, the two lists remember twice the number of pages that would fit in the cache. At any time, the policy selects a variable number of the most recent pages to exclude from the two lists. The policy adaptively decides in response to an evolving workload how many top pages from each list to maintain in the cache at any given time. It achieves such online, on-the-fly adaptation by using a learning rule that allows the policy to track a workload quickly and effectively. This allows the policy to balance between recency and frequency in an online and self-tuning fashion, in response to evolving and possibly changing access patterns. The policy is also scan-resistant. It allows one-time-only sequential read requests to pass through the cache without flushing pages that have temporal locality. The policy is extremely simple to implement and requires only constant-time overhead per request. The policy has negligible space overhead.

**42 Claims, 11 Drawing Sheets**